Sparse and Low-Rank Representation Lecture I: Motivation and Theory

Yi Ma

MSRA and **UIUC**

John Wright

Columbia University

IPAM Computer Vision Summer School, August 5, 2013

CONTEXT – Data increasingly massive, high-dimensional...



> 1B users

CONTEXT – Discovering knowledge from data



How to extract compact knowledge from such massive datasets?

CONTEXT – Good solutions impact many applications



Collaborative filtering...

Low-dimensional structures in high-dimensional data



How can we learn and exploit low-dimensional structures in high-dimensional data?

But it is not so easy...



Real application data often contain missing observations, corruptions, or even malicious errors.

Classical methods (e.g., least squares, PCA) break down...

CONTEXT - New Phenomena with High-Dimensional Data

KEY CHALLENGE: efficiently and reliably recover low-dimensional structures from high-dimensional data, despite gross observation errors.

A sobering message: human intuition is severely limited in highdimensional spaces:



Gaussian samples in 2D



As dimension grows proportionally with the number of samples...

A **new regime** of geometry, statistics, and computation...

CONTEXT - Massive High-Dimensional Data



The curse of dimensionality:

...increasingly demand inference with limited samples for very highdimensional data.

The blessing of dimensionality:

... real data highly concentrate on low-dimensional, sparse, or degenerate structures in the high-dimensional space.

But **nothing is free**: Gross errors and irrelevant measurements are now ubiquitous in massive cheap data.

Everything old ...

A long and rich history of robust estimation with error correction and missing data imputation:



R. J. Boscovich. *De calculo probailitatum que respondent diversis valoribus summe errorum post plures observationes ... , before 1756*

A. Legendre. *Nouvelles methodes pour la determination des orbites des cometes*, 1806





C. Gauss. Theory of motion of heavenly bodies, 1809

A. Beurling. Sur les integrales de Fourier absolument convergentes et leur application a une transformation functionelle, 1938



over-determined

+ dense, Gaussian



B. Logan. Properties of High-Pass Signals, 1965

underdetermined + sparse, Laplacian



... IS NEW AGAIN

Today, robust estimation in high dimensions is more urgent, more tractable, and increasingly sharply understood.

Theory – high-dimensional geometry & statistics, measure concentration, combinatorics, coding theory...

Algorithms – large scale convex optimization, parallel and distributed computing....

Applications – massive data driven methods, sensing and hashing, denoising, superresolution, MRI, bioinformatics, image classification, recognition ...





Lecture I: Motivation and Theory

- Lecture II: Data Modeling and Applications
- Lecture III: Efficient Optimization for Low-Dimensional Models
- **Tomorrow:** A Bayesian Perspective on Sparse Approximation
 - ... + Q&A and discussion

Lecture I Theory of Sparse and Low-Rank Recovery

John Wright Electrical Engineering Columbia University



Underdetermined system

$$y = Ax$$



Signal acquisition



Underdetermined system

$$y = Ax$$



Signal acquisition





Underdetermined system

$$y = Ax$$



Signal acquisition



 $y_i = \int_{\boldsymbol{u}} \boldsymbol{z}(u) \exp(-2\pi j \boldsymbol{k}(t_i)^* \boldsymbol{u}) d\boldsymbol{u}$ Observations are Fourier coefficients!



Underdetermined system

$$y = Ax$$



Signal acquisition



A few Fourier coefficients



 $oldsymbol{F}_{\Omega}$



Underdetermined system

$$y = Ax$$



Signal acquisition



[Lustig, Donoho + Pauly '10] ... brain image – Lustig '12

Underdetermined system

$$y = Ax$$



Signal acquisition



[Lustig, Donoho + Pauly '10] ... brain image – Lustig '12

Underdetermined system

$$y = Ax$$



Compression

y



Image to be compressed

Underdetermined system

$$y = Ax$$

 \approx



Compression – JPEG





(Patches of) ... input image





[Wallace '91]

Underdetermined system

$$y = Ax$$



Compression – Learned dictionary



See [Elad+Bryt '08], [Horev et. Al., '12] ... Image: [Aharon+Elad '05]

Underdetermined system $oldsymbol{y} = oldsymbol{A}oldsymbol{x}$



Recognition



Linear subspace model for images of same face under varying lighting.

[Basri+Jacobs '03], [Ramamoorthi '03], [Belhumeur+Kriegman '96]

Underdetermined system

$$y = Ax$$



Recognition



Underdetermined system

$$y = Ax$$



Recognition



 $oldsymbol{y} \in \mathbb{R}^m$ Test image



Х

Combined training dictionary



 $oldsymbol{x} \in \mathbb{R}^n$ coefficients



 $e \in \mathbb{R}^m$ corruption, occlusion

[W., Yang, Ganesh, Sastry, Ma '09]

+



One large underdetermined system: y = A'x'



Solution is **not unique** ... is there any hope?

WHAT DO WE KNOW ABOUT x?

Underdetermined system

$$y = Ax$$



Signal acquisition	Image compression	Face Recognition
i i i i i i i i i i	x^* uses just a few dictionary elements.	$\mathbf{\hat{x}^{\star} uses just a few} $ training faces. $\mathbf{\hat{z}^{\star} corrects a few} $ gross errors.

SPARSITY – More formally

A vector $x \in \mathbb{R}^n$ is **sparse** if only a few entries are nonzero:



The **number of nonzeros** is called the ℓ^0 -"norm" of x:

$$\|\boldsymbol{x}\|_0 \doteq \#\{i \mid x_i \neq 0\}.$$

SPARSITY – More formally

A vector $x \in \mathbb{R}^n$ is **sparse** if only a few entries are nonzero:



The **number of nonzeros** is called the ℓ^0 -"norm" of x:

$$\|\boldsymbol{x}\|_0 \doteq \#\{i \mid x_i \neq 0\}.$$

Geometrically

$$\|\boldsymbol{x}\|_{p} = (\sum_{i} |x_{i}|^{p})^{1/p}$$

$$\|\boldsymbol{x}\|_0 = \lim_{p \searrow 0} \|\boldsymbol{x}\|_p^p.$$



THE SPARSEST SOLUTION





Look for the sparsest \boldsymbol{x} that agrees with our observation:

minimize $\|x\|_0$ subject to Ax = y.

[Demo]

THE SPARSEST SOLUTION

Underdetermined system $oldsymbol{y} = oldsymbol{A}oldsymbol{x}$



Look for the sparsest \boldsymbol{x} that agrees with our observation:

minimize
$$\|x\|_0$$
 subject to $Ax = y$.

Theorem 1 (Gorodnitsky+Rao '97) . Suppose $y = Ax_0$, and let $k = ||x_0||_0$. If null(A) contains no 2k-sparse vectors, x_0 is the unique optimal solution to minimize $||x||_0$ subject to y = Ax.

THE SPARSEST SOLUTION

Underdetermined system $oldsymbol{y} = oldsymbol{A}oldsymbol{x}$



Look for the sparsest \boldsymbol{x} that agrees with our observation:





 $\|x\|_0$ subject to Ax = y. minimize

The cardinality $||\boldsymbol{x}||_0$ is **nonconvex**:







The cardinality $\|\boldsymbol{x}\|_0$ is **nonconvex**:

Its convex envelope* is the ℓ^1 norm: $\|\boldsymbol{x}\|_1 = \sum_i |x_i|$



* Over the set $\{\boldsymbol{x} \mid |x_i| \leq 1 \forall i\}$





The cardinality $\|\boldsymbol{x}\|_0$ is **nonconvex**:

Its convex envelope* is the ℓ^1 norm: $\|\boldsymbol{x}\|_1 = \sum_i |x_i|$






RELAX!

minimize
$$\|x\|_0$$
 subject to $Ax = y$. NP-hard, hard to appx.
[Natarjan '95],
[Amaldi+Kann '97]
minimize $\|x\|_1$ subject to $Ax = y$. Efficiently solvable

– Lecture 3!

Have we lost anything? [demo]

WHY DOES THIS WORK? Geometric intuition

minimize $\|x\|_1$ subject to Ax = y.



We see:
$$\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x} = \sum_{i \in \text{supp}(\boldsymbol{x})} \boldsymbol{a}_i x_i$$





Mutual coherence $\mu(\mathbf{A}) \doteq \max_{i \neq j} |\langle \mathbf{a}_i, \mathbf{a}_j \rangle|$

Smaller is better!



Mutual coherence

$$\mu(\boldsymbol{A}) \doteq \max_{i \neq j} |\langle \boldsymbol{a}_i, \boldsymbol{a}_j \rangle|$$

Theorem 2 (Gribonval+Nielsen '03, Donoho+Elad '03) . Suppose $y = Ax_0$ with

$$\|x_0\|_0 < \frac{1}{2}(1+1/\mu(A)).$$

Then x_0 is the unique optimal solution to

minimize $\|\boldsymbol{x}\|_1$ subject to $\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x}$.







Mutual coherence

$$\mu(\boldsymbol{A}) \doteq \max_{i \neq j} |\langle \boldsymbol{a}_i, \boldsymbol{a}_j \rangle|$$

Theorem 2 (Gribonval+Nielsen '03, Donoho+Elad '03) . Suppose $y = Ax_0$ with

$$\|x_0\|_0 < \frac{1}{2}(1+1/\mu(A)).$$

Then x_0 is the unique optimal solution to

minimize $\|\boldsymbol{x}\|_1$ subject to $\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x}$.

WHY CARE ABOUT THE THEORY?

Motivates applications

... but be careful: need to justify (and modify) the basic models [Lecture 2].

Template for stronger results

... predictions can be very sharp in high dimensions.

Generalizes to many other types of low-dimensional structure

... structured sparsity, low-rank recovery [later, Lecture 2]







MOTIVATING APPLICATIONS – Face Recognition



MOTIVATING APPLICATIONS – Face Recognition



More practicalities in Lecture 2...

WHY CARE ABOUT THE THEORY?

Motivates applications

 $\xrightarrow{x^{\star}}$

... but be careful: need to justify (and modify) the basic models [Lecture 2].

Template for stronger results

... predictions can be very sharp in high dimensions.



Generalizes to many other types of low-dimensional structure



... structured sparsity, low-rank recovery [later, Lecture 2]

LIMITATIONS OF COHERENCE?

For any
$$m \times n \ \mathbf{A} \quad \mu(\mathbf{A}) \ge \sqrt{\frac{n-m}{m(n-1)}}$$

Prev. result therefore requires

$$\|\boldsymbol{x}_0\|_0 < \frac{1}{2}(1+\mu(\boldsymbol{A})^{-1}) = O(\sqrt{m})$$



LIMITATIONS OF COHERENCE?

For any
$$m \times n \ \mathbf{A} \quad \mu(\mathbf{A}) \ge \sqrt{\frac{n-m}{m(n-1)}}$$

Prev. result therefore requires

$$\|\boldsymbol{x}_0\|_0 < \frac{1}{2}(1+\mu(\boldsymbol{A})^{-1}) = O(\sqrt{m})$$

Truth is often **much better**:







LIMITATIONS OF COHERENCE?

For any
$$m \times n \ \mathbf{A}$$
 $\mu(\mathbf{A}) \ge \sqrt{\frac{n-m}{m(n-1)}}$

Prev. result therefore requires

$$\|\boldsymbol{x}_0\|_0 < \frac{1}{2}(1+\mu(\boldsymbol{A})^{-1}) = O(\sqrt{m})$$



Plot: Fraction of correct recovery vs. fraction of nonzeros $\|\boldsymbol{x}_0\|_0/m$



STRENGTHENING THE BOUND – the RIP

Incoherence: Each pair $A_{i,j} = [a_i | a_j]$ spread.



STRENGTHENING THE BOUND – the RIP

Incoherence: Each pair $A_{i,j} = [a_i | a_j]$ spread.



Generalize to **subsets of size** *k*:

 $oldsymbol{A}_I$ well-spread (almost orthonormal) for all I of size k \Longrightarrow all k-sparse $oldsymbol{x}$, $\|oldsymbol{A}oldsymbol{x}\|_2 pprox \|oldsymbol{x}\|_2$

STRENGTHENING THE BOUND – the RIP

Incoherence: Each pair $A_{i,j} = [a_i | a_j]$ spread.

Generalize to **subsets of size** *k*:

 $oldsymbol{A}_I$ well-spread (almost orthonormal) for all I of size k \Longrightarrow all k-sparse $oldsymbol{x}$, $\|oldsymbol{A}oldsymbol{x}\|_2 pprox \|oldsymbol{x}\|_2$

A satisfies the **Restricted Isometry Property** of order k with constant δ if for all k-sparse x,

$$(1-\delta) \| \boldsymbol{x} \|_2^2 \leq \| \boldsymbol{A} \boldsymbol{x} \|_2^2 \leq (1+\delta) \| \boldsymbol{x} \|_2^2.$$



IMPLICATIONS OF RIP

Good sparse recovery

Theorem 2 (Candès+Tao '05, Candès '08) . Suppose $y = Ax_0$ with

 $\delta_{2\|\boldsymbol{x}_0\|_0} < \sqrt{2} - 1.$

Then x_0 is the unique optimal solution to

minimize $\|\boldsymbol{x}\|_1$ subject to $\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x}$.

IMPLICATIONS OF RIP

Good sparse recovery

Theorem 2 (Candès+Tao '05, Candès '08) . Suppose $y = Ax_0$ with

 $\delta_{2\|\boldsymbol{x}_0\|_0} < \sqrt{2} - 1.$

Then x_0 is the unique optimal solution to

minimize $\|\boldsymbol{x}\|_1$ subject to $\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x}$.

Again, if ... x_0 is "structured" and A is "nice" we exactly recover x_0 . Compare condition to condition $\|x_0\|_0 < \frac{1}{2}(1+\mu(A)^{-1})$

IMPLICATIONS OF RIP

Random *A* are great:

If $A \sim_{iid} \mathcal{N}(0, m^{-1/2})$ then A has RIP of order k with high probability, when $m \ge Ck \log(n/k)$.

For random $oldsymbol{A}$, ℓ^1 works even when $\|oldsymbol{x}_0\|_0 \sim m.$

Useful property for designing sampling operators *(Compressed sensing)*.

WHY CARE ABOUT THE THEORY?

Motivates applications

... but be careful: need to justify (and modify) the basic models [Lecture 2].

Template for stronger results

... predictions can be very sharp in high dimensions.

Generalizes to many other types of low-dimensional structure

... structured sparsity, low-rank recovery [later, Lecture 2]







GENERALIZATIONS – From Sparse to Low-Rank

So far: Recovering a single sparse vector:



Next: Recovering low-rank matrix (many correlated vectors):



FORMULATION – Robust PCA?



Given Y = X + E, with X low-rank, E sparse, recover X.

Numerous approaches to **robust PCA** in the literature:

- Multivariate trimming [Gnanadeskian + Kettering '72]
- Random sampling [Fischler + Bolles '81]
- Alternating minimization [*Ke* + *Kanade* '03]
- Influence functions [de la Torre + Black '03]

Can we give an efficient, provably correct algorithm?

RELATED SOLUTIONS – Matrix recovery

Classical PCA/SVD – low rank + noise [Hotelling '35, Karhunen+Loeve '72,...]

Given Y = X + Z, recover X.

Stable, efficient algorithm, theoretically optimal \rightarrow huge impact

Matrix Completion – low rank, missing data

From $Y = \mathcal{P}_{\Omega}[X]$, recover X.

[Candès + Recht '08, Candès + Tao '09, Keshevan, Oh, Montanari '09, Gross '09, Ravikumar and Wainwright '10]

Increasingly well-understood; solvable if $oldsymbol{X}$ is low rank and Ω large enough.

Our problem, with Y = X + E, looks more difficult...

WHY IS THE PROBLEM HARD?

Some very sparse matrices are also low-rank:



Can we recover X that are *incoherent* with the standard basis?

Certain sparse error patterns E make recovering X impossible:



Can we correct E whose support is not adversarial?

WHEN IS THERE HOPE? Again, (in)coherence

Can we recover X that are incoherent with the standard basis from almost all errors E?

Incoherence condition on singular vectors, singular values arbitrary:

not too cross-correlated: $\|U$

Singular vectors of
$$\boldsymbol{X}$$
 not too spiky:
$$\begin{cases} \max_{i} \|\boldsymbol{U}_{i}\|^{2} \leq \mu r/m. \\ \max_{i} \|\boldsymbol{V}_{i}\|^{2} \leq \mu r/n. \end{cases}$$
not too cross-correlated: $\|\boldsymbol{U}\boldsymbol{V}^{*}\|_{\infty} \leq \sqrt{\mu r/mn}$

Uniform model on error support, **signs and magnitudes arbitrary:**

$$\operatorname{support}(\boldsymbol{E}) \sim \operatorname{uni} \begin{pmatrix} [m] \times [n] \\ \rho mn \end{pmatrix}$$

Incoherence condition: [Candès + Recht '08]

... AND HOW SHOULD WE SOLVE IT?

Naïve optimization approach

Look for a low-rank X that agrees with the data up to some sparse error E:



... AND HOW SHOULD WE SOLVE IT?



... AND HOW SHOULD WE SOLVE IT?

Naïve optimization approach

Look for a low-rank X that agrees with the data up to some sparse error E:

min rank $(\mathbf{X}) + \gamma \|\mathbf{E}\|_0$ subj $\mathbf{X} + \mathbf{E} = \mathbf{Y}$.

Convex relaxation

Nuclear norm heuristic: [Fazel, Hindi, Boyd '01], see also [Recht, Fazel, Parillo '08]

MAIN RESULT – Correct recovery

Theorem 1 (Principal Component Pursuit). If $X_0 \in \mathbb{R}^{m \times n}$, $m \ge n$ has rank

$$r \leq \rho_r \frac{n}{\mu \log^2(m)}$$

and E_0 has Bernoulli support with error probability $\rho \leq \rho_s^{\star}$, then with very high probability

$$(X_0, E_0) = \arg \min ||X||_* + \frac{1}{\sqrt{m}} ||E||_1 \quad \text{subj} \quad X + E = X_0 + E_0,$$

and the minimizer is unique.

"Convex optimization recovers matrices of rank $O\left(\frac{n}{\log^2 m}\right)$ from errors corrupting O(mn) entries"

[Candès, Li, Ma, and W., '09].

EXAMPLE – Faces under varying illumination

58 images of one person under varying lighting:





APPLICATIONS – **Background modeling from video**

Static camera surveillance video

200 frames, 144 x 172 pixels,

Significant foreground motion









Video Y = Low-rank appx. X+ Sparse error E













BIG PICTURE – Parallelism of Sparsity and Low-Rank

	Sparse Vector	Low-Rank Matrix
Degeneracy of	individual signal	correlated signals
Measure	L ₀ norm $ x _0$	$\operatorname{rank}(X)$
Convex Surrogate	L ₁ norm $ x _1$	Nuclear norm $\ X\ _*$
Compressed Sensing	y = Ax	Y = A(X)
Error Correction	y = Ax + e	Y = A(X) + E
Domain Transform	$y \circ \tau = Ax + e$	$Y \circ \tau = A(X) + E$
Mixed Structures	Y = A(X) + B(E) + Z	

WHY CARE ABOUT THE THEORY?

Motivates applications

... but be careful: need to justify (and modify) the basic models [Lecture 2].

Template for stronger results

... predictions can be very sharp in high dimensions.

Generalizes to many other types of low-dimensional structure

... structured sparsity, low-rank recovery [later, Lecture 2]







General theory: constructing norms

Atomic norm: choose a set of atoms \mathcal{A} . Write

 $\|\boldsymbol{x}\|_{\diamond} = \inf \left\{ \sum_{i} c_{i} \mid \sum_{i} c_{i} \boldsymbol{a}_{i} = \boldsymbol{x}, \ c_{i} > 0, \boldsymbol{a}_{i} \in \mathcal{A} \right\}$

[Chandrasekharan et. al. '12]

General theory: constructing norms

Atomic norm: choose a set of atoms \mathcal{A} . Write

$$\|\boldsymbol{x}\|_{\diamond} = \inf \left\{ \sum_{i} c_{i} \mid \sum_{i} c_{i} \boldsymbol{a}_{i} = \boldsymbol{x}, \ c_{i} > 0, \boldsymbol{a}_{i} \in \mathcal{A} \right\}$$

E.g., sparsity $\mathcal{A} = \{ e_i \mid i = 1 \dots n \}, \| x \|_{\diamond} = \| x \|_{\ell^1}$ low-rank $\mathcal{A} = \{ uv^* \mid \| u \|_2 = \| v \|_2 = 1 \}, \| x \|_{\diamond} = \| x \|_*$

[Chandrasekharan et. al. '12]
General theory: constructing norms

Atomic norm: choose a set of atoms \mathcal{A} . Write

$$\|\boldsymbol{x}\|_{\diamond} = \inf \left\{ \sum_{i} c_{i} \mid \sum_{i} c_{i} \boldsymbol{a}_{i} = \boldsymbol{x}, \ c_{i} > 0, \boldsymbol{a}_{i} \in \mathcal{A} \right\}$$

E.g., sparsity $A = \{ e_i \mid i = 1 ... n \}, \|x\|_{\diamond} = \|x\|_{\ell^1}$ low-rank $\mathcal{A} = \{ uv^* \mid ||u||_2 = ||v||_2 = 1 \}$, $||x||_\diamond = ||x||_*$

column sparsity



 $\mathcal{A} = \{ \boldsymbol{u}\boldsymbol{e}_i^* \mid \|\boldsymbol{u}\|_2 = 1, \ i = 1 \dots n \}$ e.g., [Xu+Caramanis+Sanghavi'12]

sinusoids $\bigwedge \bigwedge \bigwedge \qquad \mathcal{A} = \{e^{2\pi f t + \xi} \mid f \in [0, 1], \ \xi \in [0, 2\pi)\}$ [Tang + Recht '12] [Candes + Fernandez-Garza '12]

General theory: constructing norms

Atomic norm: choose a set of atoms \mathcal{A} . Write

$$\|\boldsymbol{x}\|_{\diamond} = \inf \left\{ \sum_{i} c_{i} \mid \sum_{i} c_{i} \boldsymbol{a}_{i} = \boldsymbol{x}, \ c_{i} > 0, \boldsymbol{a}_{i} \in \mathcal{A} \right\}$$

E.g., sparsity $\mathcal{A} = \{ e_i \mid i = 1 \dots n \}, \| \boldsymbol{x} \|_{\diamond} = \| \boldsymbol{x} \|_{\ell^1}$ low-rank $\mathcal{A} = \{ \boldsymbol{u} \boldsymbol{v}^* \mid \| \boldsymbol{u} \|_2 = \| \boldsymbol{v} \|_2 = 1 \}, \| \boldsymbol{x} \|_{\diamond} = \| \boldsymbol{x} \|_*$



Observe: $\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x}_0$ with $\boldsymbol{A} \sim \mathcal{N}(0,1)$ random. When does $\min \|\boldsymbol{x}\|_{\diamond}$ s.t. $\boldsymbol{A}\boldsymbol{x} = \boldsymbol{y}$

uniquely recover x_0 ?



Observe: $y = Ax_0$ with $A \sim \mathcal{N}(0, 1)$ random. When does $\min ||x||_\diamond$ s.t. Ax = y uniquely recover x_0 ?

Observe: $\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x}_0$ with $\boldsymbol{A} \sim \mathcal{N}(0,1)$ random. When does $\min \|\boldsymbol{x}\|_\diamond$ s.t. $\boldsymbol{A}\boldsymbol{x} = \boldsymbol{y}$

uniquely recover x_0 ?



Recovery iff the **descent cone** $D(\|\cdot\|_{\diamond}, \boldsymbol{x}_{0}) = \{\boldsymbol{v} \mid \|\boldsymbol{x}_{0} + t\boldsymbol{v}\|_{\diamond} \leq \|\boldsymbol{x}_{0}\|_{\diamond} \text{ for some } t > 0\}$ has $D(\|\cdot\|_{\diamond}, \boldsymbol{x}_{0}) \cap \text{null}(\boldsymbol{A}) = \{\boldsymbol{0}\}.$

More likely if descent cone is "small". Can we make this precise?

Observe: $y = Ax_0$ with $A \sim \mathcal{N}(0,1)$ random. When does $\min ||x||_\diamond$ s.t. Ax = y uniquely recover x_0 ?

The **statistical dimension** of a cone *C* is

 $\delta(C) = \mathbb{E}_{\boldsymbol{g} \sim \mathcal{N}(0,1)} \left[\left\| P_C \boldsymbol{g} \right\|^2 \right].$



Many nice properties. E.g., if *C* a subspace, $\delta(C) = \dim(C)$.

Observe: $\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x}_0$ with $\boldsymbol{A} \sim \mathcal{N}(0,1)$ random. When does $\min \|\boldsymbol{x}\|_\diamond$ s.t. $\boldsymbol{A}\boldsymbol{x} = \boldsymbol{y}$ uniquely recover \boldsymbol{x}_0 ?

The statistical dimension of a cone *C* is $\delta(C) = \mathbb{E}_{\boldsymbol{g} \sim \mathcal{N}(0,1)} \left[\|P_C \boldsymbol{g}\|^2 \right].$



Many nice properties. E.g., if *C* a subspace, $\delta(C) = \dim(C)$.

Sharp **phase transition** at $m = \delta(C)$:

$$m > \delta(C) \implies \mathbb{P}[\text{recovery}] > 1 - \exp\left(-c(m - \delta(C))^2/n\right)$$
$$m < \delta(C) \implies \mathbb{P}[\text{recovery}] < \exp\left(-c(m - \delta(C))^2/n\right)$$

[Amelunxen, McCoy, Lotz, Tropp '13]

General theory: decomposing two structures

Observe: $y = x_0 + z_0$ with regularizers $||x||_{\diamond,1}$, $||z||_{\diamond,2}$. Does

 $\min \|\boldsymbol{x}\|_{\diamond,1} + \|\boldsymbol{z}\|_{\diamond,2} \quad \text{s.t.} \quad \boldsymbol{x} + \boldsymbol{z} = \boldsymbol{y}$

uniquely recover $oldsymbol{x}_0$, $oldsymbol{z}_0$?

Variant: $\min \|\boldsymbol{x}\|_{\diamond,1}$ s.t. $\|\boldsymbol{z}\|_{\diamond,2} \leq 1, \ \boldsymbol{x} + \boldsymbol{z} = \boldsymbol{y}$

$$C_{2} = D(\|\cdot\|_{\diamond,2}, \boldsymbol{z}_{0})$$

$$\boldsymbol{0}$$

$$C_{1} = D(\|\cdot\|_{\diamond,1}, \boldsymbol{x}_{0})$$

General theory: decomposing two structures

Observe: $y = x_0 + z_0$ with regularizers $||x||_{\diamond,1}$, $||z||_{\diamond,2}$. Does

 $\min \|\boldsymbol{x}\|_{\diamond,1} + \|\boldsymbol{z}\|_{\diamond,2} \quad \text{s.t.} \quad \boldsymbol{x} + \boldsymbol{z} = \boldsymbol{y} \qquad C_2 = D(\|\cdot\|_{\diamond,2}, \boldsymbol{z}_0)$ uniquely recover $\boldsymbol{x}_0, \, \boldsymbol{z}_0$? Variant: $\min \|\boldsymbol{x}\|_{\diamond,1} \quad \text{s.t.} \quad \|\boldsymbol{z}\|_{\diamond,2} \le 1, \, \boldsymbol{x} + \boldsymbol{z} = \boldsymbol{y} \qquad C_1 = D(\|\cdot\|_{\diamond,1}, \boldsymbol{x}_0)$

In a random incoherence model (C_2 randomly rotated), **phase transition** at

$$\delta(C_1) + \delta(C_2) = n$$

 $n > \delta(C_1) + \delta(C_2) \implies \mathbb{P}[\text{recovery}] > 1 - \exp\left(-c(n - \delta(C_1) - \delta(C_2))^2/n\right)$ $n < \delta(C_1) + \delta(C_2) \implies \mathbb{P}[\text{recovery}] < \exp\left(-c(n - \delta(C_1) - \delta(C_2))^2/n\right)$

[Amelunxen, McCoy, Lotz, Tropp '13]

General theory: statistical estimation

Observe: noisy measurements $y = Ax_0 + z$. Noise-aware optimization:

$$\min \|oldsymbol{x}\|_{\diamond} + rac{\gamma}{2} \|oldsymbol{A}oldsymbol{x} - oldsymbol{y}\|_2^2$$

E.g., Basis pursuit denoising: min $\|x\|_1 + \frac{\lambda}{2} \|Ax - y\|_2^2$ Noise-aware RPCA: min $\|L\|_* + \lambda \|S\|_1 + \frac{\gamma}{2} \|L + S - D\|_F^2$

When does $\min \|\boldsymbol{x}\|_{\diamond} + \frac{\gamma}{2} \|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{y}\|_2^2$ produce $\hat{\boldsymbol{x}} \approx \boldsymbol{x}_0$?

General theory for **decomposable regularizers** $\|\cdot\|_{\diamond}$

[Negahbhan, Agarwal, Yu, Wainwright '12]

A suite of models and theoretical guarantees

For robust recovery of a family of low-dimensional structures:

- [Zhou et. al. '09] Spatially contiguous sparse errors via MRF
- [Bach '10] structured relaxations from **submodular functions**
- [Negahban+Yu+Wainwright '10] **geometric analysis** of recovery
- [Becker+Candès+Grant '10] algorithmic templates
- [Xu+Caramanis+Sanghavi '11] column sparse errors L_{2,1} norm
- [Recht+Parillo+Chandrasekaran+Wilsky '11] compressive sensing of various structures
- [Candes+Recht '11] compressive sensing of decomposable structures

$$X^0 = \arg \min ||X||_\diamond$$
 s.t. $\mathcal{P}_Q(X) = \mathcal{P}_Q(X^0)$

- [McCoy+Tropp'11] decomposition of sparse and low-rank structures $(X_1^0, X_2^0) = \arg \min ||X_1||_{(1)} + \lambda ||X_2||_{(2)}$ s.t. $X_1 + X_2 = X_1^0 + X_2^0$
- [W.+Ganesh+Min+Ma, I&I'13] superposition of decomposable structures

 $(X_1^0,\ldots,X_k^0) = \arg\min\sum\lambda_i ||X_i||_{(i)}$ s.t. $\mathcal{P}_Q(\sum_i X_i) = \mathcal{P}_Q(\sum_i X_i^0)$

Take home message: Let the data and application tell you the structure...

Thank you! Question time ...

Next: **Modeling** low-dimensional structure in **real data**.

Later: Solving the **optimization** problems efficiently.

Tomorrow: A **Bayesian** way of approaching sparse modeling.